

# Exploiting FAIR Data to Enhance Data Analysis

3rd International Network Meeting of EUSMAT  
29.06.2022

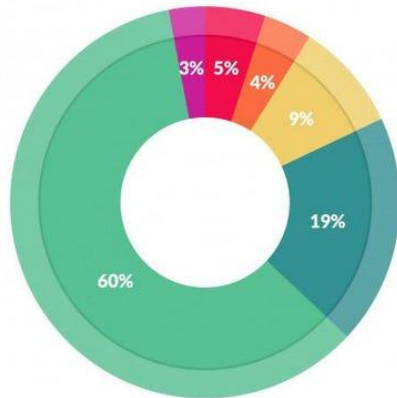
Dr. Marius Politze (RWTH Aachen University)



This work is licensed under a [Creative Commons Attribution 4.0 International License](https://creativecommons.org/licenses/by/4.0/).

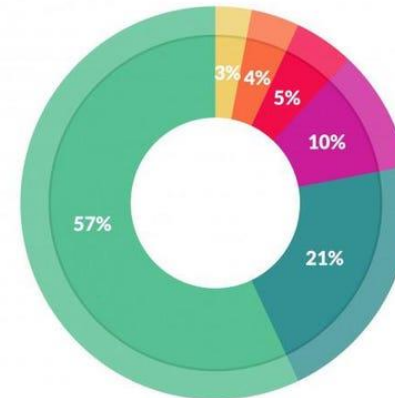


# Why Should I care?



What data scientists spend the most time doing

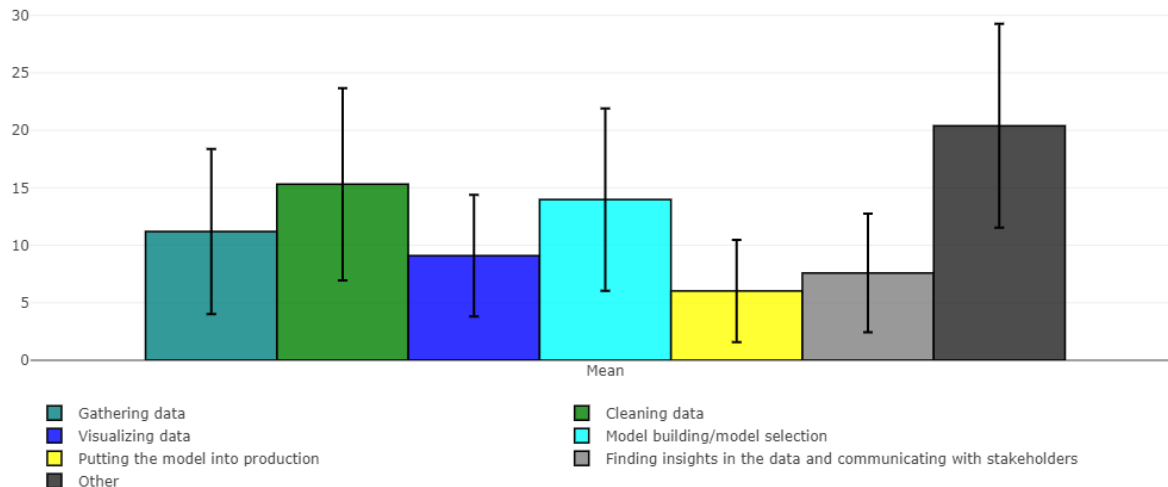
- Building training sets: 3%
- Cleaning and organizing data: 60%
- Collecting data sets: 19%
- Mining data for patterns: 9%
- Refining algorithms: 4%
- Other: 5%



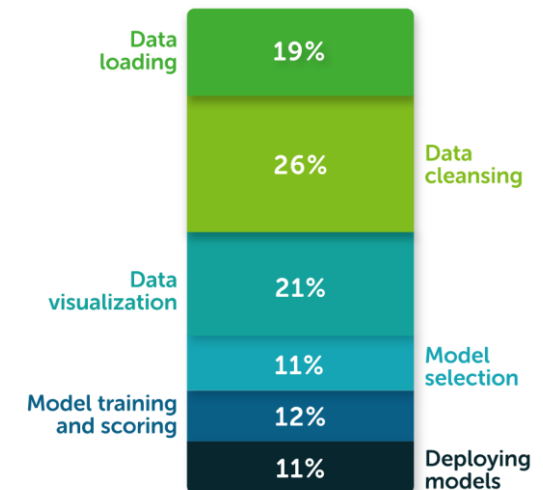
What's the least enjoyable part of data science?

- Building training sets: 10%
- Cleaning and organizing data: 57%
- Collecting data sets: 21%
- Mining data for patterns: 3%
- Refining algorithms: 4%
- Other: 5%

Press, G (2016): Cleaning Big Data: Most Time-Consuming, Least Enjoyable Data Science Task, Survey Says, <https://www.forbes.com/sites/gilpress/2016/03/23/data-preparation-most-time-consuming-least-enjoyable-data-science-task-survey-says/#1b963786f637>



Mooney, P (2018): Kaggle Machine Learning & Data Science Survey, <https://www.kaggle.com/code/paultimothymooney/2018-kaggle-machine-learning-data-science-survey/notebook>



Anaconda Inc. (2020): 2020 State of Data Science, <https://www.anaconda.com/state-of-data-science-2020>

Data preparation and cleansing takes valuable time away from real data science work and has a negative impact on overall job satisfaction.

Anaconda Inc. (2020): 2020 State of Data Science,  
<https://www.anaconda.com/state-of-data-science-2020>

... and (most of the time) this is because most data collection are a mess!

Marius Politze, today

# FAIR Principles<sup>1</sup>

- Findable:  
Persistent identifiers; indexed and searchable metadata
- Accessible:  
Retrievable using standard protocols; “tombstones”
- Interoperable:  
Vocabularies for data and metadata; qualified references to other (meta)data
- Reusable:  
Licenses; provenance; meets community standards



FAIR guiding principles for data resources, Sangya Pundir, CC-BY-SA-4.0

<sup>1</sup>Wilkinson, M. D. et. al. (2016). The FAIR Guiding Principles for scientific data management and stewardship. doi:10.1038/sdata.2016.18

# Recap FAIR Principles (1)

---

## To be Findable:

- F1. (Meta)data are assigned a globally unique and persistent identifier
- F2. Data are described with rich metadata
- F3. Metadata clearly and explicitly include the identifier of the data they describe
- F4. (Meta)data are registered or indexed in a searchable resource

## To be Accessible:

- A1. (Meta)data are retrievable by their identifier using a standardised communications protocol
  - A1.1 The protocol is open, free, and universally implementable
  - A1.2 The protocol allows for an authentication and authorisation procedure, where necessary
- A2. Metadata are accessible, even when the data are no longer available

## Recap FAIR Principles (2)

---

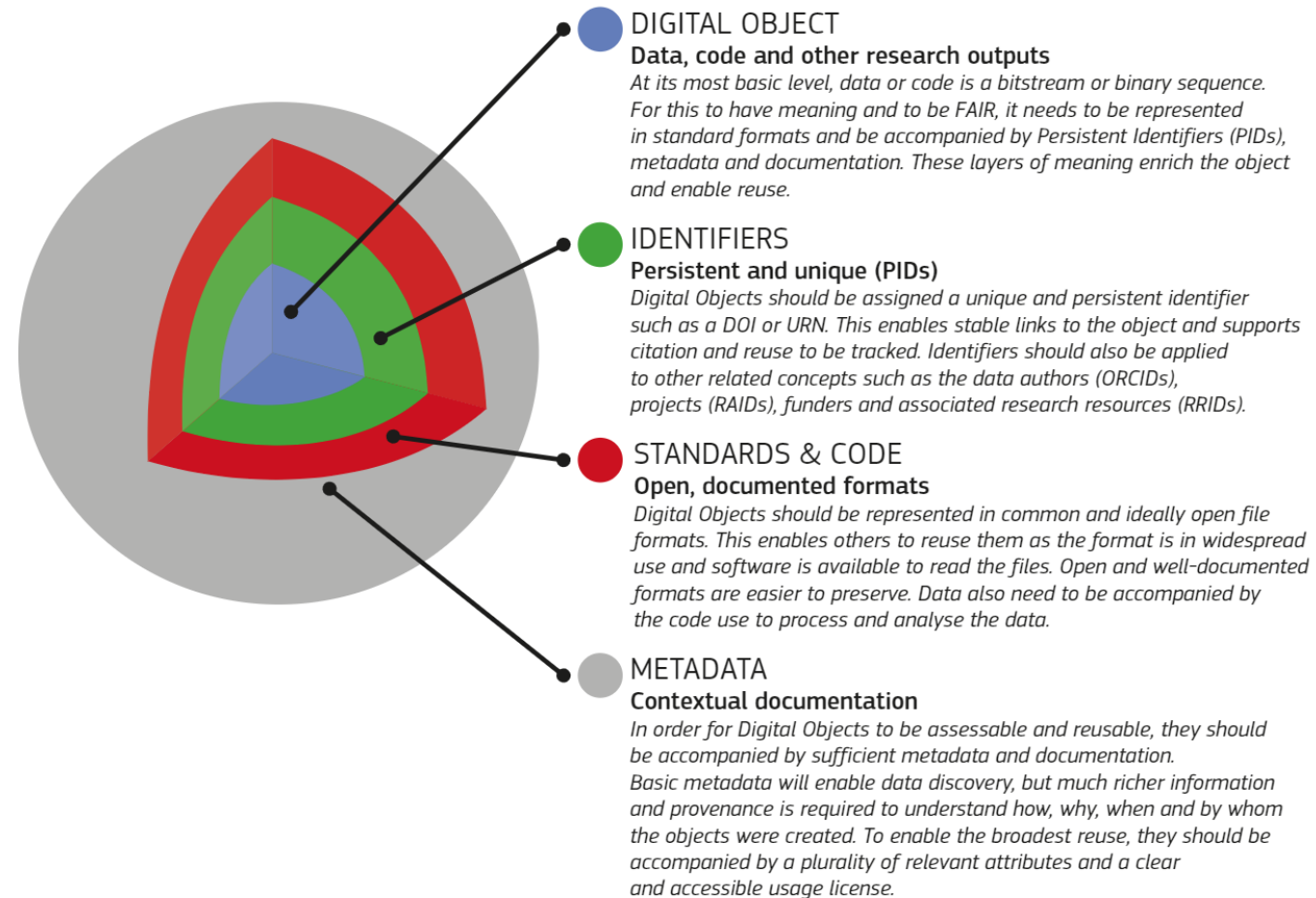
### To be Interoperable:

- I1. (Meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation.
- I2. (Meta)data use vocabularies that follow FAIR principles
- I3. (Meta)data include qualified references to other (meta)data

### To be Re-usable:

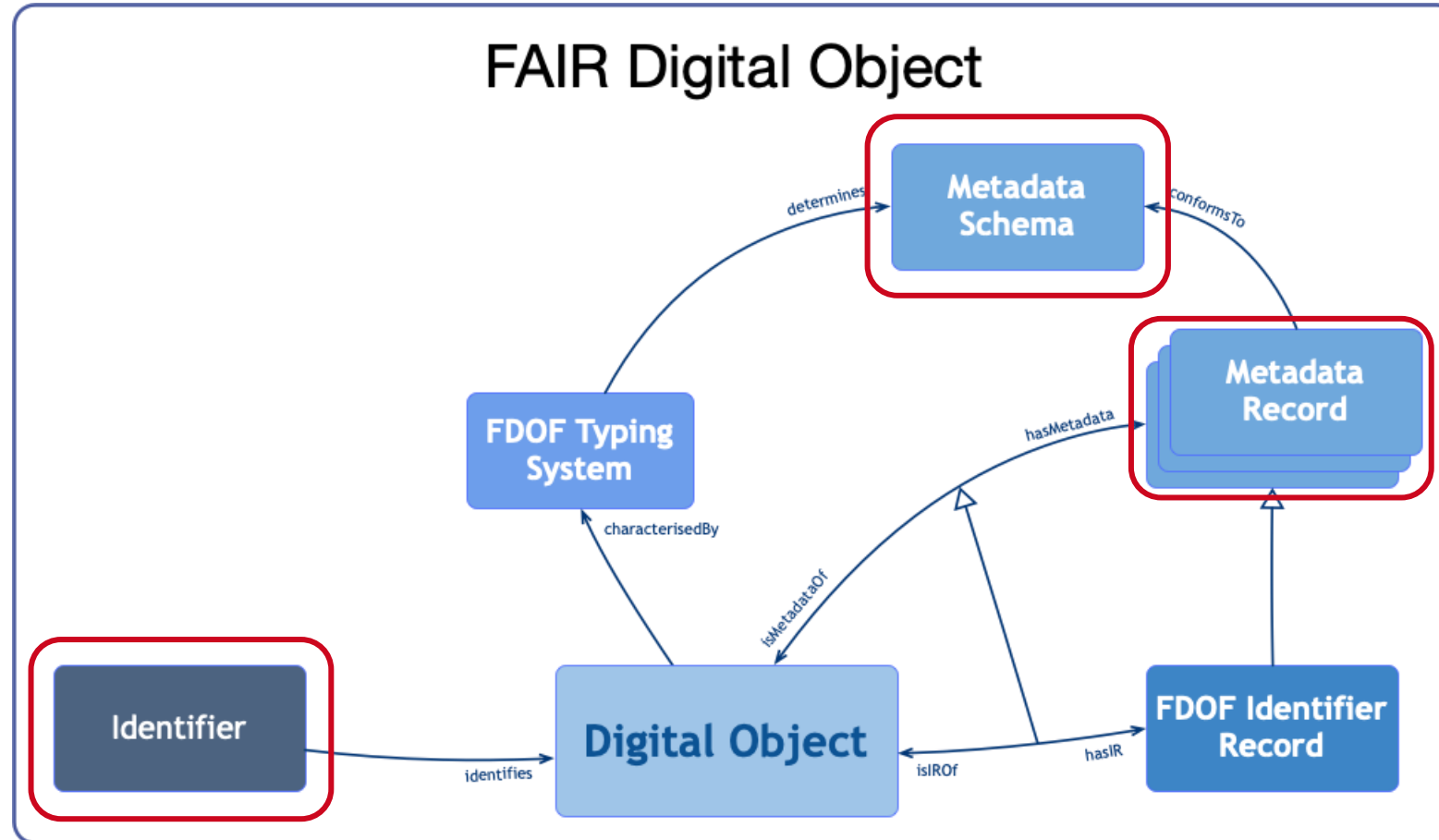
- R1. (Meta)data are richly described with a plurality of accurate and relevant attributes
  - R1.1. (Meta)data are released with a clear and accessible data usage license
  - R1.2. (Meta)data are associated with detailed provenance
  - R1.3. (Meta)data meet domain-relevant community standards

# Implementing FAIR Principles: FAIR Digital Objects<sup>1</sup>



<sup>1</sup> European Commission, Directorate-General for Research and Innovation, Turning FAIR into reality : final report and action plan from the European Commission expert group on FAIR data, Publications Office, 2018, doi:10.2777/1524

# FAIR Digital Object – A Structured Set of Links



Bonino da Silva Santos, L. O. (2021): FAIR Digital Object Framework Documentation. <https://fairdigitalobjectframework.org/>



# Persistent Identifiers - Eternally and Globally Unique

---

- What is eternity?
  - Whatever the community defines it to be
- What is globally unique?
- Examples
  - **DOI – Published Documents (and Data)**
  - **ORCID – Researchers**
  - ROR – Research Organizations
  - **ePIC/Handle – (Unpublished) Data, real world things**
  - pURL – Mostly digital concepts
  - W3ID – Mostly digital concepts
  - ...

*Alice: How long is forever?*  
*White Rabbit: Sometimes, just one second.*  
– Alice in Wonderland<sup>1</sup>

<sup>1</sup> All credits for finding this quote go to York Sure-Vetter (Director of NFDI e.V.)

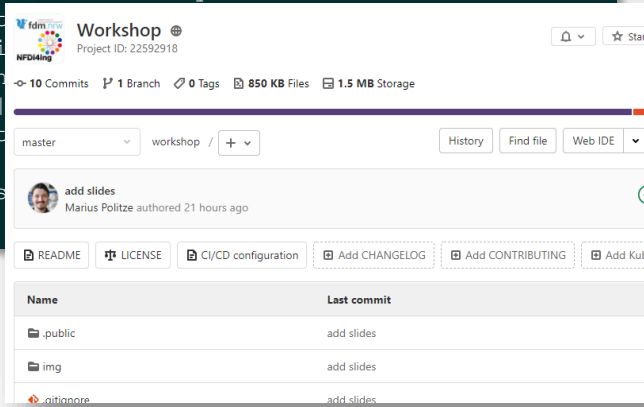
# (Structured) Metadata Records

- Text-Files
- Metadata In Files
- Databases
- XML / JSON
- Need for standardization across different storage methods

```

---
sensor: sensorA
deviceUnderTest: Prototype C
title: 'This is the title: it contains a colon'
author:
- name: Author One
  affiliation: University of Somewhere
- name: Author Two
  affiliation: University of Somewhere
tags: [nothing, something]
abstract: |
  This is the abstract.

It consists of several slides.
---
    
```

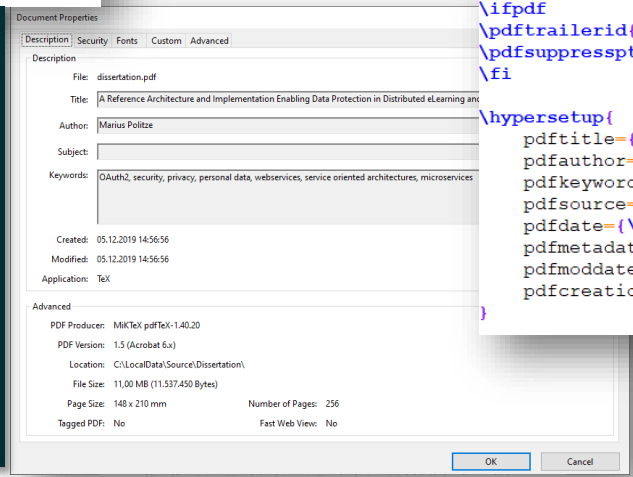
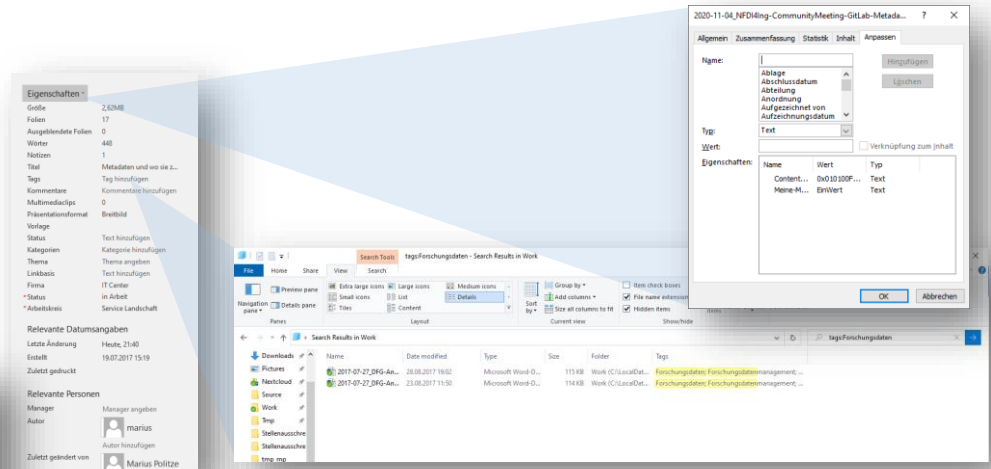


```

year=2019&month=12&day=24.js
year=2019&month=12.js
year=2019.js
year=2020&month=04&day=30.js
year=2020&month=04.js
year=2020.js
years.js
    
```

```

Experiment1
├── README.md
├── .index.yml
├── run1
│   ├── .index.yml
│   ├── sensorA
│   │   ├── .index.yml
│   │   ├── dataset1.hdf5
│   │   ├── dataset1.json
│   │   ├── dataset2.csv
│   │   └── sensorB
│   │       ├── .index.yml
│   │       └── dataset.abc
└──
    
```



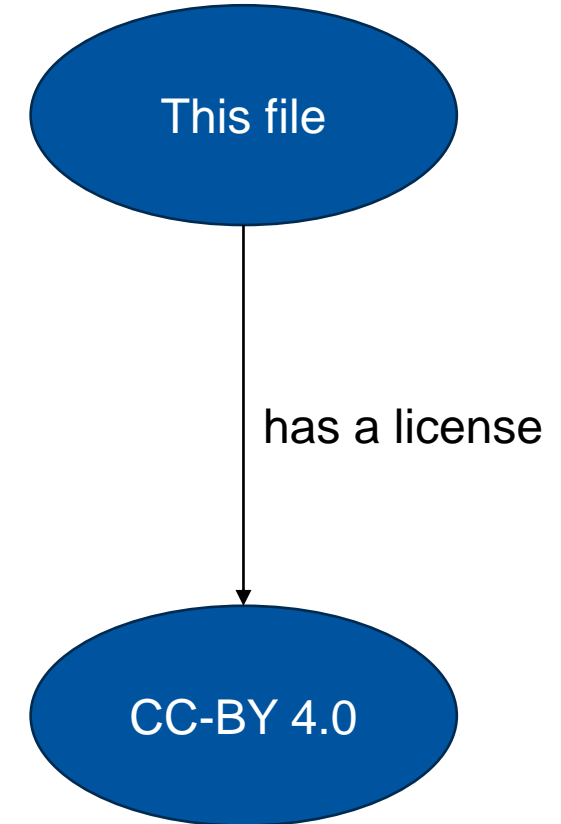
```

% % % % % % % % % % PDF Metadaten % % % % % % % % % %
\ifpdf
\pdftrailerid(\gitHash)
\pdfsuppressptexinfo--1
\fi

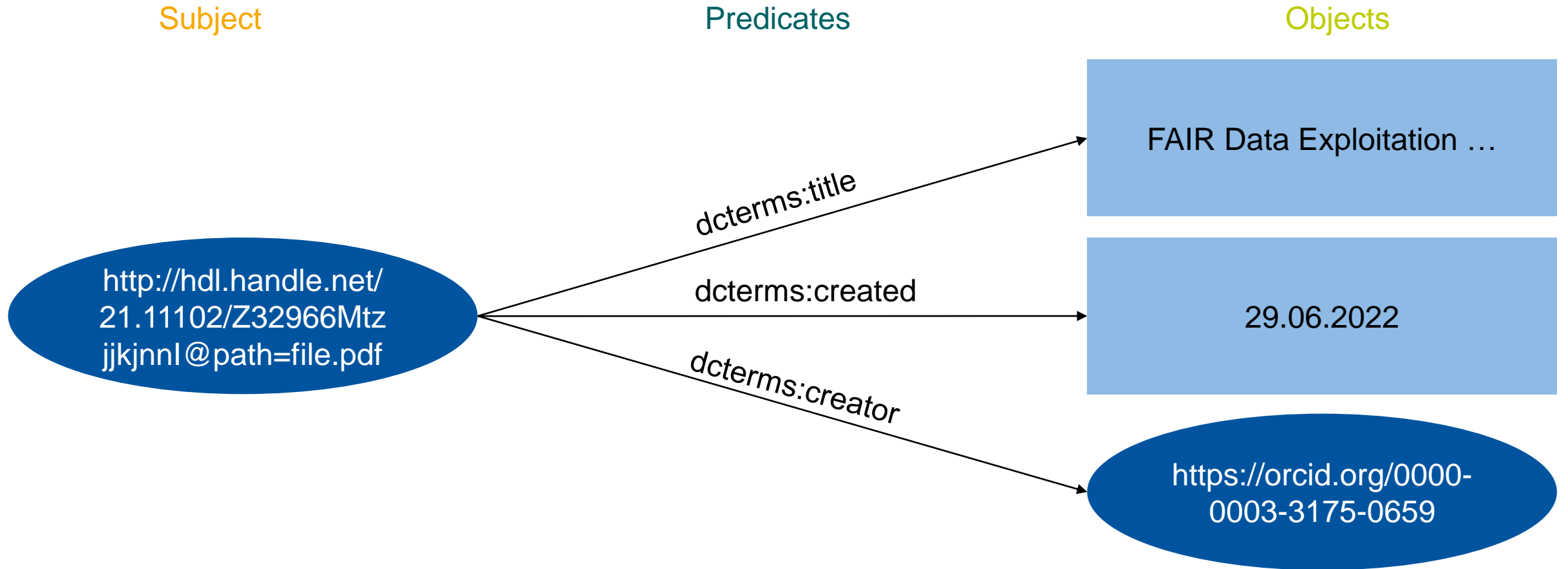
\hypersetup{
pdftitle={\thetitle},
pdfauthor={\theauthorfirstname\ theauthorlastname},
pdfkeywords={\thekeywords},
pdfsource={\gitAbbrevHash},
pdfdate={\gitStrictIsoDate},
pdfmetadate={\gitStrictIsoDate},
pdfmoddate={\gitPdfDate},
pdfcreationdate={\gitPdfDate}
}
    
```

# Semi Structured Data on the Web: Linked Data and the Semantic Web

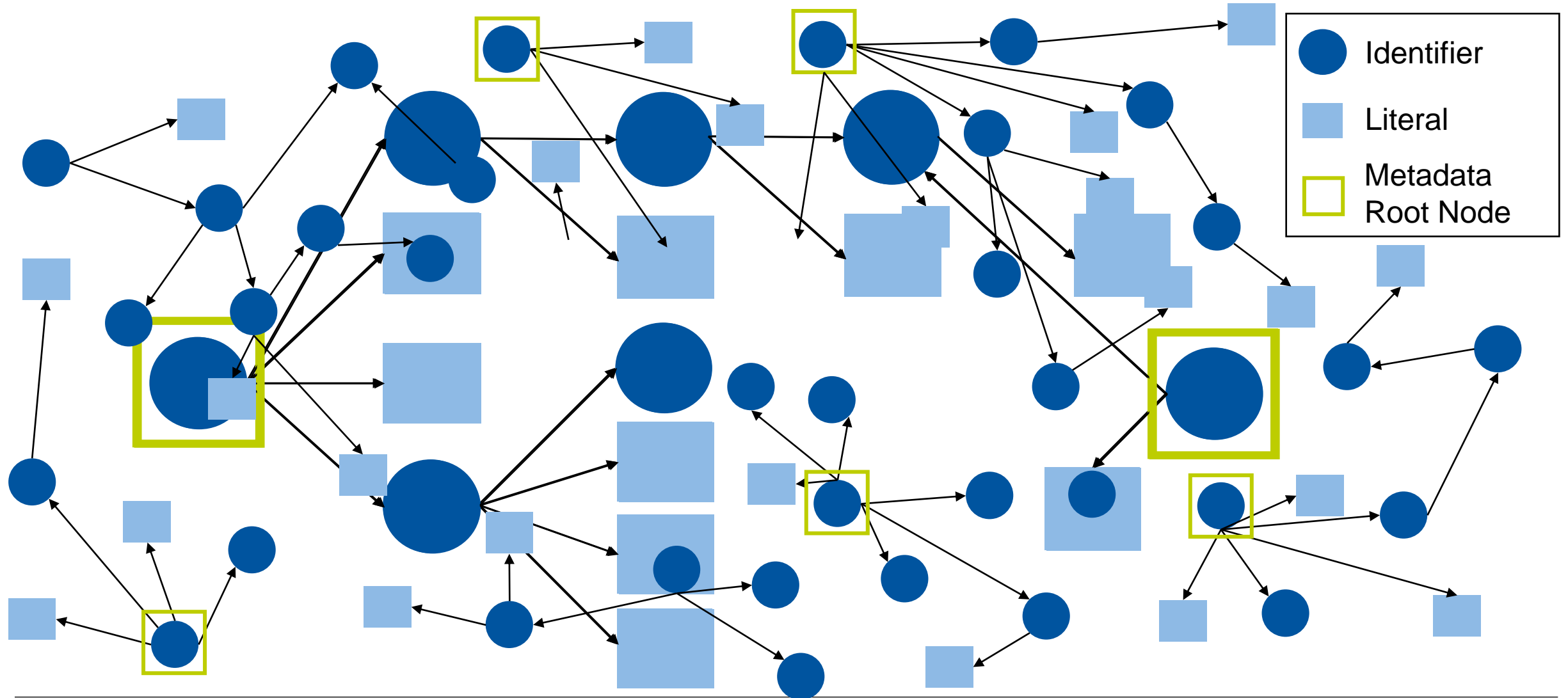
- Structured because the framework is defined → RDF
- RDF defines (meta)data to be saved in triples of “subject predicate object”
  - Think of very simple sentences of the form “A has property of value B
- Unstructured because the terms used can (and should) be freely defined
- Linked Data Terms to uniquely identify concepts
  - Real world things
  - Digital Objects
  - Relations between these things
  - Abstract properties
- Open Specification allows adding to various formats
  - HTML
  - PDF
  - HDF5
  - ...



# RDF Creates a Metadata Graph

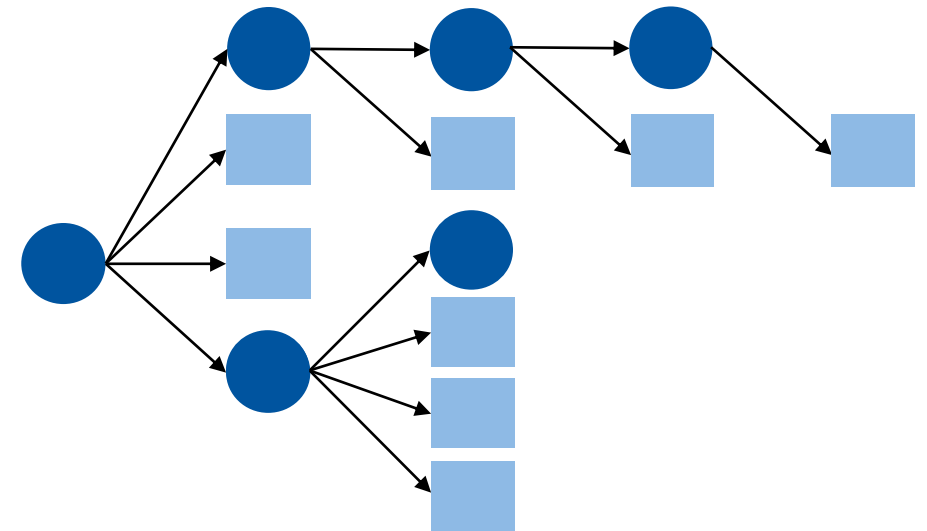


# Metadata as a Knowledge Graph in the Semantic Web



## Bringing Order to the Chaos

- Metadata Schema
  - “Open World Definition”: Everything is valid if it not explicitly forbidden
  - Set of information to describe things in a given context
  - What do predicates (relations) mean? What are their semantics?
- Metadata Standard
  - A Metadata Schema that is endorsed by some standardization body
- Application Profiles
  - “Cloesed World Definition”:
  - May consist of Information from multiple schemas and standards
  - Define your data structures based on your field of application

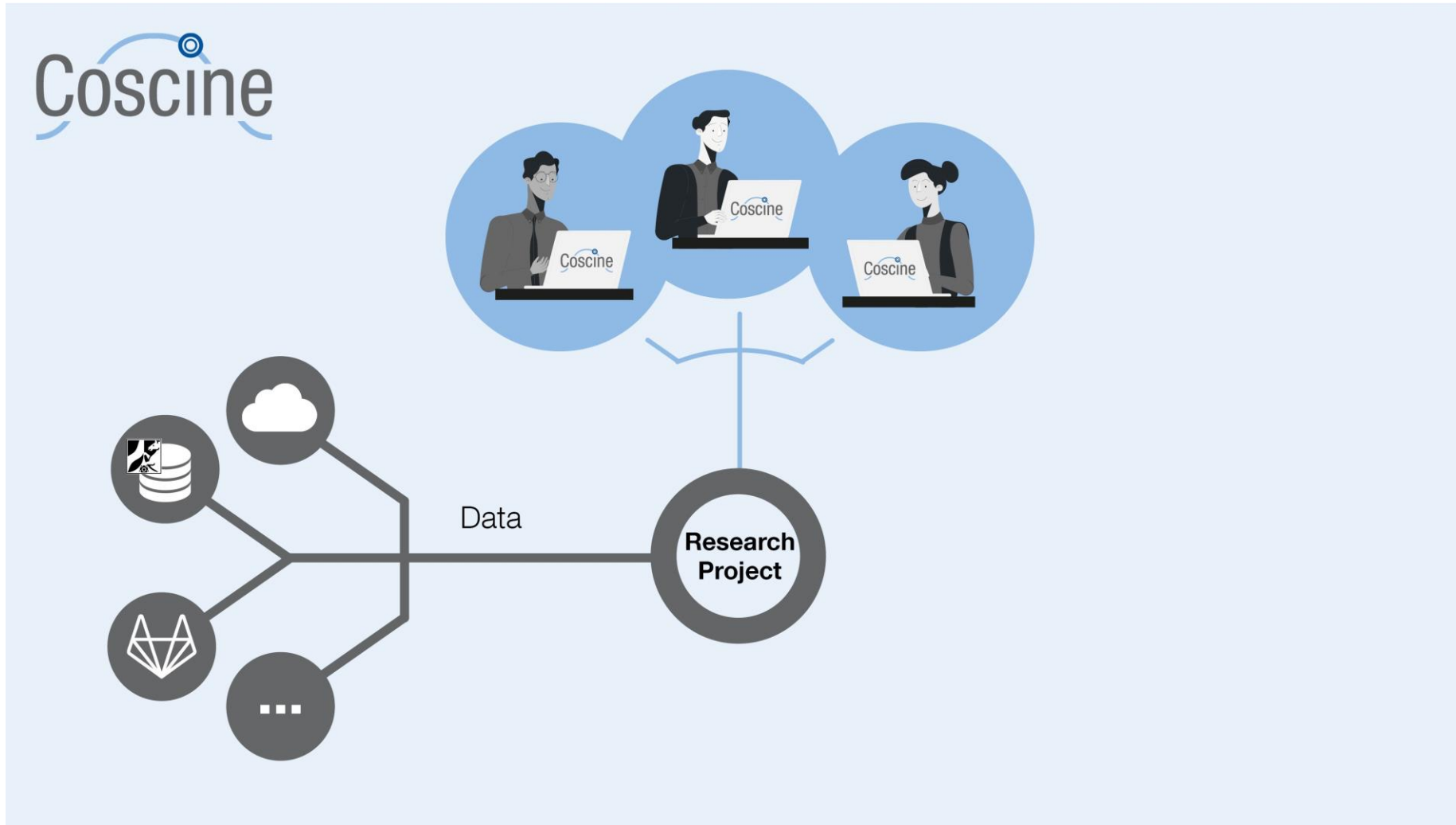


## Ok, nice but...

---

... is there an app for that?

# Coscine - Short Introduction

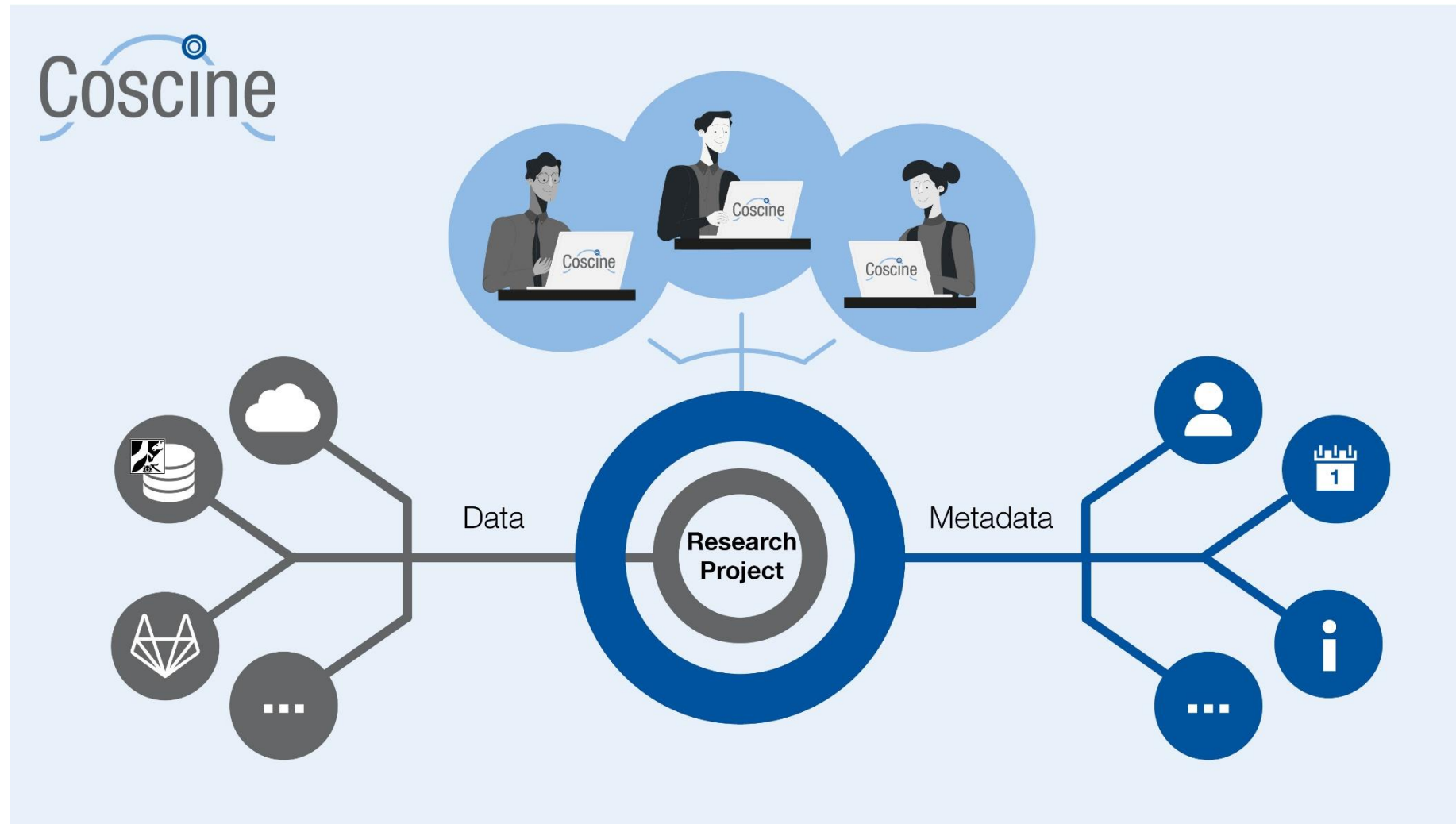


## Coscine...

- ...maps a project structure
- ...integrates different storage systems
- ...maintains authorizations at the project level



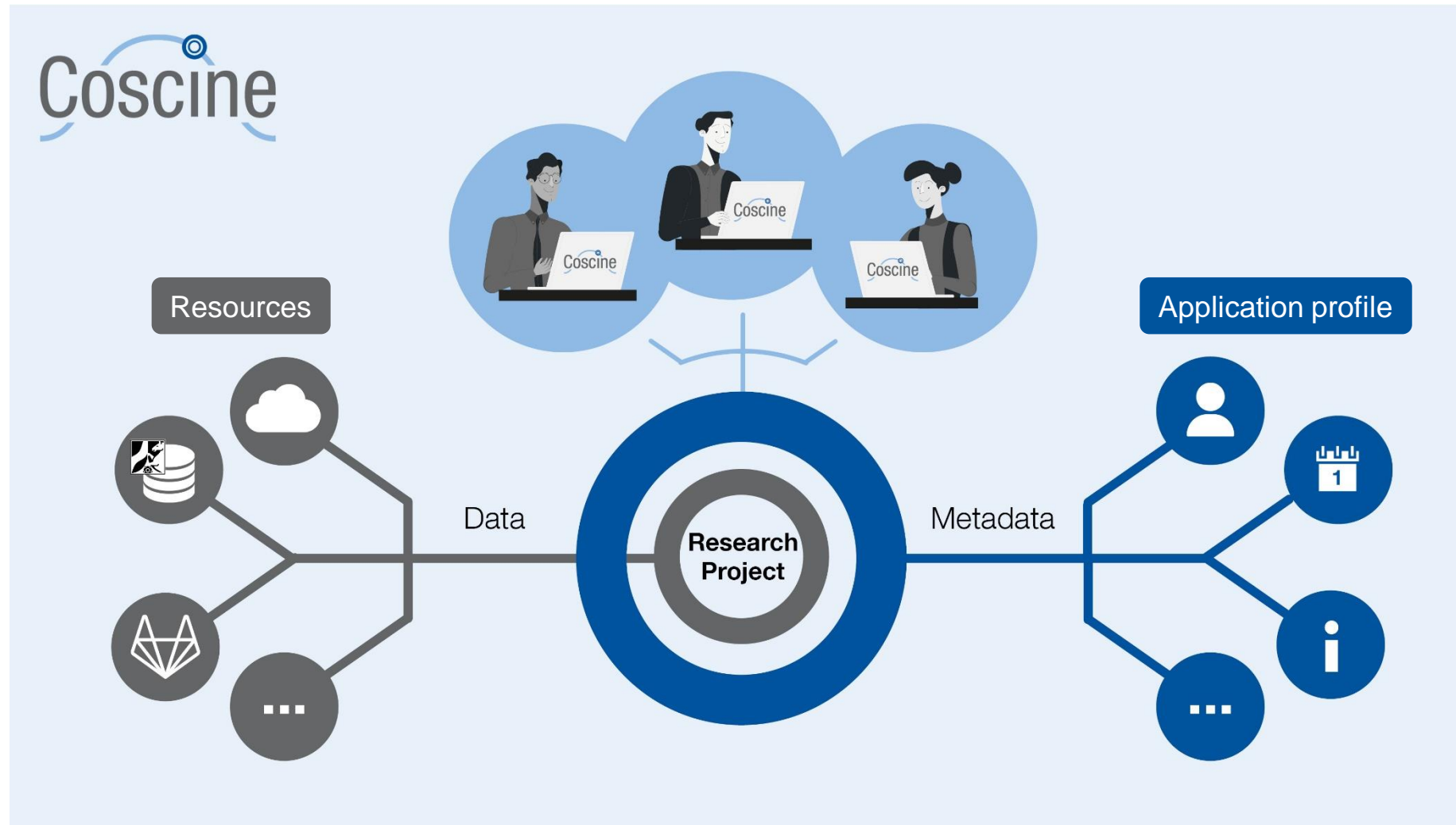
# Coscine - Short Introduction



## Coscine...

- ...maps a project structure
- ...integrates different storage systems
- ...maintains authorizations at the project level
- ...helps to describe all data with structured metadata

# Coscine - Short Introduction



## Coscine...

- ...maps a project structure
- ...integrates different storage systems
- ...maintains authorizations at the project level
- ...helps to describe all data with structured metadata

# Coscine – Metadata Management

## 1. Project Level

<b>Project Name: *</b>	Autonomous Driving in the City Center	✓
<b>Display Name: *</b>	Autonomous driving - CC	✓
<b>Project Description: *</b>	In this project, autonomous driving in the city center is recorded using measurement data from ten vehicles. The vehicles drove on a test route under changing weather conditions and high pedestrian and traffic volumes.	✓
<hr/>		
<b>Project Metadata</b>		
<b>Principal Investigators (PIs): *</b>	Conny Taylor	✓
<b>Project Start: *</b>	<input type="calendar"/> Tuesday, February 8, 2022	
<b>Project End: *</b>	<input type="calendar"/> Wednesday, April 17, 2024	
<b>Discipline: *</b>	Electrical Engineering and Information Technology 408	▼
<b>Participating Organizations: *</b>	RWTH Aachen University <input type="checkbox"/> TU Dortmund University <input type="checkbox"/>	▼
<b>Project Keywords:</b>	autonomous driving <input type="checkbox"/> vehicles <input type="checkbox"/>	▼
<b>Grant ID:</b>	DFG_007	

## 2. Resource Level

Options:

- **Pre-implemented** application profiles (e.g. EngMeta, see figure)
- Create and request an **individual** application profile

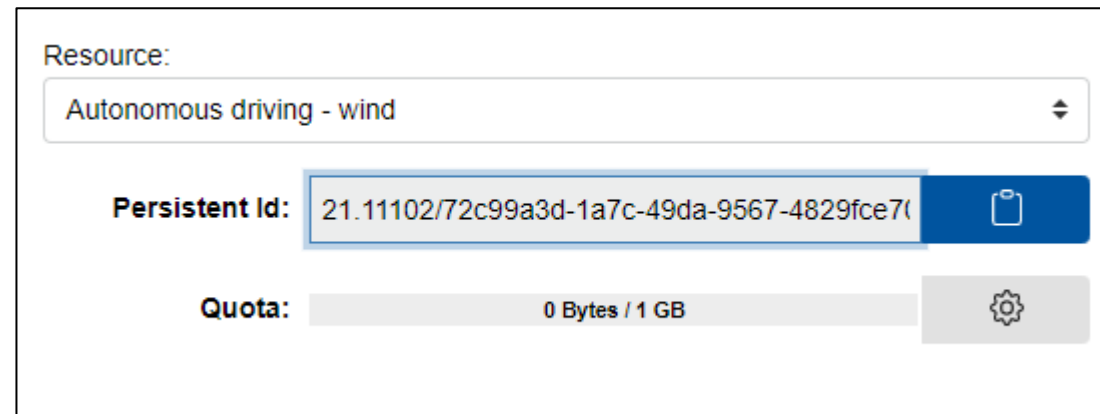
The screenshot displays a metadata management form with the following fields and values:

Field	Value	Lock	Visibility	Plus
Application Profiles *	Engmeta	Yes	Yes	Yes
Contact *	Conny Research	Yes	Yes	Yes
Creator *	Conny Research	Yes	Yes	Yes
Worked	Yes	Yes	Yes	Yes
Worked Note	Vehicle passed the crossroad without a crash	Yes	Yes	Yes
Title *	Autonomous driving -windy and rain	Yes	Yes	Yes
Type	Dataset	Yes	Yes	Yes
Keywords	autonomous driving	Yes	Yes	Yes
Subject Area	Traffic and Transport Systems, Logistics, Intelligent and Automated Traffic	Yes	Yes	Yes
Creation Date *	Tuesday, February 8, 2022	Yes	Yes	Yes
Publication Date *	Friday, February 11, 2022	Yes	Yes	Yes
Embargo End Date *	Thursday, April 25, 2024	Yes	Yes	Yes
Version *	1	Yes	Yes	Yes

## Coscine – Persistent Identifiers (PIDs)


---


- Coscine uses PIDs to uniquely reference resources
- Each resource in a project is automatically assigned a PID
- The URL contains the handle-prefix followed by a PID
  - Example: <http://hdl.handle.net/21.11102/7599d318-99f3-4385-ace9-7aeb9cf3bXXX>
- PIDs can be used to link resources and make them accessible to others



Resource:

Autonomous driving - wind

**Persistent Id:** 21.11102/72c99a3d-1a7c-49da-9567-4829fce7( 

**Quota:** 0 Bytes / 1 GB 

# The AIMS Project - Creating Your Own Application Profile

<https://coscine.rwth-aachen.de/coscine/apps/aimsfrontend/>


The screenshot shows the AIMS application profile editor interface. The browser address bar displays the URL: `coscine.rwth-aachen.de/coscine/apps/aimsfrontend/?token=eyJhbGciOiJIUzI1NiIsInR5cCI6IkpXVCJ9.eyJvc2VzYySWQzOTg0OTgwNmI2Ni02YjJLTQyM2M0GE4ZC1mNDkwYWQ2MjM...`. The interface is divided into four main panels:

- Available Application Profiles:** A search bar with the text "Search Application Profiles" and a list of profile URLs, including `https://purl.org/coscine/ap/din4000-128/`, `https://purl.org/coscine/ap/sfb985/arc`, `https://purl.org/coscine/ap/cwd/`, `https://purl.org/coscine/ap/engmeta/`, `https://purl.org/coscine/ap/ltt/`, `https://purl.org/coscine/ap/radar/`, `https://purl.org/coscine/ap/sfb1394/to`, `https://purl.org/coscine/ap/sfb1394/At`, and `https://purl.org/coscine/ap/sfb1394/At`.
- Vocabulary Terms:** A search bar with the text "resolution" and a list of terms with green plus icons: `OEO_00000514`, `dcattemporalResolutio`, `OEO_00000516`, `OEO_00000515`, `dcatspatialResolutionI`, and `hasAssumption`.
- Detail View:** A table showing the configuration for the "radar application profile". The table has two columns: the property name and a status indicator (a red 'x').

Creator	x
Title	x
Production Date	x
Subject Area	x
Resource	x
Rights	x
Rightsholder	x
uses instrument	x
- Field Properties:** A form for configuring the field properties. It includes fields for "Term IRI" (with a link icon), "Field Name", "Minimum Required Entries", "Maximum Possible Entries", "Position On Metadata Form", "Property Type", and "Datatype".

# Vielen Dank für Ihre Aufmerksamkeit

Dr. Marius Politze

 0000-0003-3175-0659

[politze@itc.rwth-aachen.de](mailto:politze@itc.rwth-aachen.de)



This work is licensed under a [Creative Commons Attribution 4.0 International License](https://creativecommons.org/licenses/by/4.0/).

